

Generator liczb losowych

Rachunek Prawdopodobieństwa i Statystyka

Paweł Pyciński

Uniwersytet Jagielloński

TABLE OF CONTENTS

- 1 Wprowadzenie
- 2 Sposoby generowania liczb pseudolosowych
- 3 Własny generator
- 4 Własny generator - kod źródłowy
- 5 Modyfikacje generatora dla uzyskania zadanych rozkładów
- 6 Test poprawności generatora
- 7 Sources

Cel Projektu

Celem projektu jest stworzenie generatora całkowitych liczb pseudolosowych o rozkładzie równomiernym. Na podstawie stworzonego generatora należy stworzyć generatory o rozkładzie jednostajnym (na przedziale $[0,1]$), Bernoulliego, dwumianowego, Poissona, wykładniczego i normalnego. Następnie należy przetestować powstałe generatory.

Definicja

Generator liczb pseudolosowych – program lub podprogram, który na podstawie niewielkiej ilości informacji generuje deterministycznie ciąg bitów, który pod pewnymi względami jest nieodróżnialny od ciągu uzyskanego z prawdziwie losowego źródła.

Generator liczb pseudolosowych nie bez powodu jest **pseudolosowy**, problem z otrzymaniem liczb losowych wynika z deterministycznego charakteru komputera i wykonywanych przez niego operacji. Gdy człowiek dokonuje rzutu kością, nie wie co wypadnie. Taka sama operacja na komputerze wymaga działania, którego wynik jest nieprzewidywalny – żadna z operacji wykonywanych przez procesor nie posiada takiej cechy.

Problem starano się rozwiązać wykorzystując zewnętrzne źródła sygnałów losowych (np. generatory białego szumu), jednakże w tego typu urządzenia nie są standardowo wyposażano komputery osobiste. Próbowano także wykorzystać szumy kart dźwiękowych, jednakże system ten nie rozpowszechnił się z prostej przyczyny – różne karty dźwiękowe szumią różnie, a te z górnej półki nie szumią prawie wcale.

Sposoby generowania liczb pseudolosowych

Jest wiele sposobów generowania liczb pseudolosowych. Jedną z grup generatorów są generatory liniowe. tworzą ciąg liczb według schematu:

$$X_{n+1} = (a_1X_n + a_2X_{n-1} + \dots + a_kX_{n-k+1} + c) \bmod(m)$$

gdzie a_1, \dots, a_k, c, m -parametry generatora (ustalone liczby)

Generatory używające operacji modulo nazywamy **kongruencyjnymi**. Każdy kolejny wyraz (liczba pseudolosowa) w generatorze liniowym to suma pewnych poprzednich wyrazów pomnożonych każdy z każdą o jakiś skalar i brane z nich jest modulo.

Generator moltiplikatywny tworzy liczby według schematu:

$$X_{i+1} = (aX_i + c) \bmod(m) \iff c = 0$$

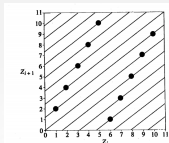
Kolejny wyraz tworzymy po przez pomnożenie poprzedniego przez jakiś skalar. Gdy $c \neq 0$ to generator jest kongruentnie mieszany.

Własny generator

Swój generator postanowiłem zbudować na bazie generatora moltiplikatywnego. Jest to jeden z łatwiejszych generatorów, prosty do implementacji.

Posiada on niestety dwie poważne wady:

1. Generator generuje liczby ciągu w sposób deterministyczny przez co łatwo jest wyliczyć kolejną liczbę.
2. Wybierając złe czynniki możemy spowodować, że okres generatora będzie mały przez co będzie działał niepoprawnie lub będzie generował bardzo mało liczb losowych.
3. Generowane liczby lokalizują się na hiperpłaszczyznach, których położenie uzależnione jest od parametrów generatora.



Przez wyżej wymienione czynniki nie może być on stosowany w kryptografii.

Przed zaimplementowaniem pozostał jeszcze wybór m oraz a dla naszego generatora.

Niech $m = 2^{32}$, jest to liczba o 1 większa od zakresu double'a, dzięki czemu nasza kongruencja potencjalnie będzie mogła zwracać wszystkie liczby które jesteśmy w stanie zapisać na 4 bajtach float'a w większości języków programowania. Ponadto niech $a = 747796405$.

TABLE 5. LCGs with Good Figures of Merit, for $m = 2^e$ and $c = 0$

m	a, a^*	$M_8(m, a)$	$M_{16}(m, a)$	$M_{32}(m, a)$
2^{30}	177911525, 17372909	0.74878 *	0.53850	0.53850
	156051869, 52274357	0.69501	0.67940 *	0.64413
	143133861, 233896749	0.69305	0.66791	0.66791 *
2^{31}	594156893, 452271861	0.75913 *	0.50244	0.50244
	558177141, 413965533	0.68978	0.68749 *	0.59450
	602169653, 448899357	0.67295	0.67116	0.67116 *
2^{32}	741103597, 887987685	0.75652 *	0.53707	0.53707
	1597334677, 851723965	0.70068	0.67686 *	0.64694
	747796405, 204209821	0.66893	0.66001	0.66001 *

Własny generator - kod źródłowy

```

1  class generator:
2
3      def __init__(self, seed):
4          self.value = seed
5          self.a = 747796405
6          self.m = 4294967296
7
8      def generateRandom(self):
9          self.value = (self.a*self.value) % self.m
10         return self.value
11

```

Listing 1: Klasa generatora

Klasa generatora posiada konstruktor który jako argument przyjmuje ziarno czyli dowolną liczbę początkową która rozpocznie budowanie pseudolosowy ciąg. Jest także metoda która zwraca kolejną wygenerowaną liczbę.

Aby uzyskać liczby z rozkładu jednostajnego na przedziale $[0, 1]$ wystarczy podzielić przez ustalone wcześniej $m = 4294967296$, zauważmy że po takiej operacji liczby będą należały to przedziału $[0, 1]$.

```
1  def uniformDistribution(self):  
2      return self.generateRandom()/self.m  
3
```

Listing 2: Metoda rozkładu jednostajnego

Rozkład Bernoulliego, jest rozkładem dwupunktowym, aby uzyskać ten rozkład skorzystam z metody którą przygotowałem dla rozkładu jednostajnego. Ustalmy dowolne $P \in [0, 1]$. Jeśli wylosowana liczba przez metodę rozkładu jednostajnego będzie większa od p to zwrócimy 0, w przeciwnym razie 1.

```
1  def bernoulliDistribution(self, probability):  
2      rand = self.uniformDistribution()  
3      if ( rand >= probability):  
4          return 0  
5      else:  
6          return 1  
7
```

Listing 3: Metoda rozkładu Bernoulliego

Rozkład dwumianowy jest to liczba sukcesów w n próbach Bernoulliego. W implementacji wykorzystałem wcześniej przygotowaną metodę generowania próby Bernoulliego, wywołanie jej *samples* razy daje nam rozkład Dwumianowy

```
1  def binomialDistribution(self, probablity, n):  
2      counter = 0  
3      for i in range(n):  
4          counter += self.bernoulliDistribution(probablity)  
5      return counter  
6
```

Listing 4: Metoda rozkładu Dwumianowego

Rozkład Poissona modeluje zdarzenia rzadkie. Jest on parametryzowany zmienną λ która jest równa oczekiwanej liczbie zdarzeń w danym przedziale czasu. Jest wiele algorytmów generujących ten rozkład na potrzeby naszego generatora wystarczy zastosować najprostszy z nich czyli **Algorytm Knutha**

```
1  def poissonDistribution(self, lambdapoiss):
2      limit = math.exp(-lambdapoiss)
3      n = 0
4      p = self.uniformDistribution()
5      while(p>=limit):
6          n+=1
7          p*=self.uniformDistribution()
8      return n
9
```

Listing 5: Metoda rozkładu Poissona

Rozkład wykładniczy modeluje czas między kolejnymi zdarzeniami, jeśli w jednostce czasu zachodzi średnio λ niezależnych zdarzeń.

Metodę generującą rozkład wykładniczy możemy uzyskać stosując metodę odwórcanej dystrybuanty. Dystrybuenta określonego rozkładu prawdopodobieństwa jest funkcją $F : \mathbb{R} \rightarrow \mathbb{R}$ niemalejąca i prawostronnie ciągła. Dystrybuenta jednoznacznie definiuje rozkład prawdopodobieństwa i ma następujący związek z gęstością prawdopodobieństwa: $F(x) = \int_x^{-\infty} f(y) dy$. Jeśli uda się znaleźć F^{-1} to $U = F(x) \rightarrow x = F^{-1}(U)$ zmienna losowa x ma rozkład o dystrybuancie F , U jest zmienną losową o rozkładzie jednostajnym. Krótki dowód dlaczego tak jest:

Niech $X = F^{-1}(U)$ zmienna losowa

$$\begin{aligned} P\{X \leq x\} &= P\{F^{-1}(U) \leq x\} \\ &= P\{U \leq F(x)\} \\ &= F(x) \end{aligned}$$

Przejdźmy teraz do rozkładu wykładniczego. Jego gęstość prawdopodobieństwa dana jest wzorem:

$$f(x) = e^{-x}, x \in [0, \infty)$$

Natomiast dystrybuanta jest całką z funkcji gęstości.

$$F(x) = \int_x^0 e^{-x} dx = 1 - e^{-x}$$

$$F(x) = 1 - e^{-x} = U$$

$$e^{-x} = 1 - U$$

$$F^{-1}(x) = x = -\ln(1 - U)$$

$$U \in (0, 1) \rightarrow x \in (0, \infty)$$

```
1 def exponentialDistribution(self):  
2     return -math.log(1-self.uniformDistribution())  
3
```

Listing 6: Metoda rozkładu Poissona

Rozkład normalny jest jednym z najważniejszych rozkładów prawdopodobieństwa, odgrywający ważną rolę w statystyce. Przyczyną jego znaczenia jest częstość występowania w naturze. Jeśli jakaś wielkość jest sumą lub średnią bardzo wielu drobnych losowych czynników, to niezależnie od rozkładu każdego z tych czynników jej rozkład będzie zbliżony do normalnego (na podstawie CTG).

Jest wiele algorytmów aby uzyskać rozkład normalny. Ja w swojej implementacji zastosowałem polarny algorytm Boxa-Mullera nazwany inaczej sposobem polarnym. Polega on na wylosowaniu dwóch zmiennych (x, y) z przedziału $(-1, 1)$ tak aby $0 < x^2 + y^2 < 1$ a następnie należy podstawić do wzoru:

$$x\sqrt{\frac{-2 \ln s}{s}} \quad \text{lub} \quad y\sqrt{\frac{-2 \ln s}{s}}$$

wzory te stosujemy na zmianę dlatego przyda się drobna modyfikacja obecnego generatora o dodanie nowej zmiennej którą będziemy zmieniać w zależności o zastosowanego wzoru.

```
1  def normalDistribution(self):
2  if(self.whichOne == 1):
3      self.whichOne = 0
4      return self.prevValue
5  else:
6      x = self.uniformDistribution()*2-1
7      y=self.uniformDistribution()*2-1
8      s=x*x+y*y
9      while (s>=1 or s==0):
10         x = self.uniformDistribution()*2-1
11         y=self.uniformDistribution()*2-1
12         s=x*x+y*y
13     s=math.sqrt((math.log(s)*(-2))/s)
14     self.prevValue=y*s
15     self.whichOne=1
16 return x*s
17
```

Listing 7: Metoda rozkładu normalnego

Kolejnym etapem jest przetestowanie napisanego wcześniej generatora oraz sprawdzenie czy otrzymane wyniki są zgodne z oczekiwanymi. Do tego celu wykorzystam test χ^2

Definicja

Test chi-kwadrat – każdy test statystyczny, w którym statystyka testowa ma rozkład chi kwadrat, jeśli teoretyczna zależność jest prawdziwa. Test chi-kwadrat służy sprawdzaniu hipotez. Innymi słowy wartość testu oceniana jest za pomocą rozkładu chi kwadrat. Test najczęściej wykorzystywany w praktyce. Można go wykorzystywać do badania zgodności zarówno cech mierzalnych, jak i niemierzalnych.

Teza zerowa: Otrzymany rozkład jest pożądanym rozkładem

Teza alternatywna: Otrzymany rozkład nie jest oczekiwanym rozkładem.

Należy ukształtować dane w taki sposób aby każda próbka miała co najmniej 5 elementów, przyjmijmy także powszechnie stosowany 5% stopień akceptacji.

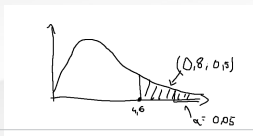
df \ p	0,99	0,95	0,90	0,80	0,50	0,30	0,20	0,10	0,05	0,02	0,01
1	0,000	0,004	0,016	0,064	0,455	1,074	1,642	2,706	3,841	5,412	6,635
2	0,020	0,103	0,211	0,446	1,386	2,408	3,219	4,605	5,991	7,824	9,210
3	0,115	0,352	0,584	1,005	2,366	3,665	4,642	6,251	7,815	9,837	11,34
4	0,297	0,711	1,064	1,649	3,357	4,878	5,989	7,779	9,488	11,67	13,28
5	0,554	1,145	1,610	2,343	4,351	6,064	7,289	9,236	11,07	13,39	15,09
6	0,872	1,635	2,204	3,070	5,348	7,231	8,558	10,64	12,59	15,03	16,81
7	1,239	2,167	2,833	3,822	6,346	8,383	9,803	12,02	14,07	16,62	18,47
8	1,646	2,733	3,490	4,594	7,344	9,524	11,03	13,36	15,51	18,17	20,09
9	2,088	3,325	4,168	5,380	8,343	10,66	12,24	14,68	16,92	19,68	21,67
10	2,558	3,940	4,865	6,179	9,342	11,78	13,44	15,99	18,31	21,16	23,21

Figure: tabela rozkładu chi kwadrat

expected	quantity	chi kwadrat
142	146	0.11267605633802817
142	146	0.11267605633802817
142	127	1.5845070422535212
142	156	1.380281690140845
142	132	0.704225352112676
142	141	0.007042253521126761
142	152	0.704225352112676
suma testu chi kwadrat = 4.605633802816902		

Figure: Rezultaty przeprowadzonego testu przy 6 stopniach swobody

nasza suma testu $\chi^2 = 4.605633802816902$ z tabelki możemy odczytać, że znajdujemy się w przedziale $[0.8, 0.5]$ zatem przy 5% stopniu akceptacji nie możemy odrzucić tezy zerowej. Nie jest to być może wynik bardzo zadawalający lecz wciąż odcinamy dość duże pole pod wykresem funkcji testu.



```
expected    quantity
600          606
400          394
suma testu chi kwadrat= 0.15
```

Figure: Rezultaty przeprowadzonego testu przy 1 stopniu swobody

nasza suma testu $\chi^2 = 0.15$ z tabelki możemy odczytać, że znajdujemy się w przedziale $[0.8, 0.5]$ zatem przy 5% stopniu akceptacji nie możemy odrzucić tezy zerowej.

```
expected      quantity
9765           9903
97656          97152
439453         438760
1171875        1170964
2050781        2049968
2460937        2461908
2050781        2050934
1171875        1171843
439453         440917
97656          97878
9765           9773
suma testu Chi kwadrat= 12.458520109015799
```

Figure: Rezultaty przeprowadzonego testu przy 10 stopniach swobody

nasza suma testu $\chi^2 = 12.458520109015799$ z tabelki możemy odczytać, że znajdujemy się w przedziale $[0.2, 0.3]$ zatem przy 5% stopniu akceptacji nie możemy odrzucić tezy zerowej.

```

67.3794699085467      50
336.89734995427335    334
842.2433748856834     824
1403.7389581428056    1407
1754.673697678507    1782
1754.673697678507    1706
1462.2280813987559    1471
1044.44862957054      1083
652.7803934815875     679
362.6557741564375     352
181.3278870782187     183
82.42176685373579     78
34.34240285572325     25
13.208616482970477    17
4.717363029632314     6
suma testu Chi kwadrat= 13.759110251832714
    
```

Figure: Rezultaty przeprowadzonego testu przy 14 stopniach swobody

nasza suma testu $\chi^2 = 13.759110251832714$ z tabelki możemy odczytać, że znajdujemy się w przedziale $[0.8, 0.5]$ zatem przy 5% stopniu akceptacji nie możemy odrzucić tezy zerowej.

```
expected      quantity
86.10666495797781  88
77.91253239626417  105
70.49817464607867  71
63.78938632300601  71
57.71902361860703  56
52.226332302616974  60
47.25633967418796  48
42.759304376622566  51
38.69021856915683  50
35.008357473362835  32
31.676871785877523  22
28.662418878189545  38
25.934829092406215  34
23.466803793176698  25
21.23364215377446  26
19.212993941920804  19
17.38463583114813  17
15.730268998951525  12
14.23333598602239  15
12.878854983630811  10
11.653269890648035  19
10.544314639530178  7
9.540890433391237  14
8.632954665513726  7
7.811420409564923  9
7.06806547458412  8
6.395450114531812  7
5.786842568810742  8
5.236151688543286  9
suma testu chi kwadrat 40.03574762305712
```

Figure: Rezultaty przeprowadzonego testu przy 29 stopniach swobody

nasza suma testu $\chi^2 = 40.03574762305712$ z tabelki możemy odczytać, że znajdujemy się w przedziale $[0.1, 0.05)$ zatem przy 5% stopniu akceptacji nie możemy odrzucić tezy zerowej.

130.000000000079	150
209.999999999991	205
349.9999999999477	381
570.000000000039	586
890.000000000019	887
1310.0	1344
1889.999999999916	1872
2600.000000000023	2574
3429.999999999995	3524
4359.99999999997	4360
5320.000000000003	5315
6240.000000000001	6203
7030.000000000003	7050
7609.999999999945	7671
7930.000000000004	7762
7930.000000000004	7858
7609.999999999945	7548
7030.000000000003	7042
6240.000000000001	6116
5320.000000000003	5329
4359.99999999997	4396
3429.999999999995	3544
2600.000000000023	2698
1889.999999999916	1853
1310.0	1313
890.000000000019	896
570.000000000039	545
349.9999999999477	349
209.999999999991	215
130.000000000079	151
suma testu chi kwadrat wynosi 31.44101220198305	

Figure: Rezultaty przeprowadzonego testu przy 29 stopniach swobody

nasza suma testu $\chi^2 = 31.44101220198305$ z tabelki możemy odczytać, że znajdujemy się w przedziale $[0.5, 0.3]$ zatem przy 5% stopniu akceptacji nie możemy odrzucić tezy zerowej.

Sources

- http://home.agh.edu.pl/~chwiej/mn/generatory_16.pdf
- https://pl.wikipedia.org/wiki/Generator_liczb_pseudolosowych
- https://eduinf.waw.pl/inf/alg/001_search/0022.php
- <https://www.ams.org/journals/mcom/1999-68-225/S0025-5718-99-00996-5/S0025-5718-99-00996-5.pdf>
- <http://staff.iiar.pwr.wroc.pl/grzegorz.mzyk/kmi/kmi03.pdf>
- <https://math.stackexchange.com/questions/785188/simple-algorithm-for-generating-poisson-distribution/785200>
- https://pl.wikipedia.org/wiki/Rozk\T1\lad_wyk\T1\ladniczy
- https://pl.xcv.wiki/wiki/Marsaglia_polar_method
- https://www.naukowiec.org/tablice/statystyka/rozklad-chi-kwadrat_247.html
- wykład Rachunek Prawdopodobieństwa i Statystyka Prof. Adam Roman